# PRICE INTELLIGENCE (PI) – DATA GATHERING BY UTILIZING WEB CRAWLING

**Big Data Applications and Utilising Non-Traditional Data Sources and Methods for Official Statistics**

**10 June 2021**

**Ms. Mazliana Mustapa**
**Department of Statistics Malaysia**
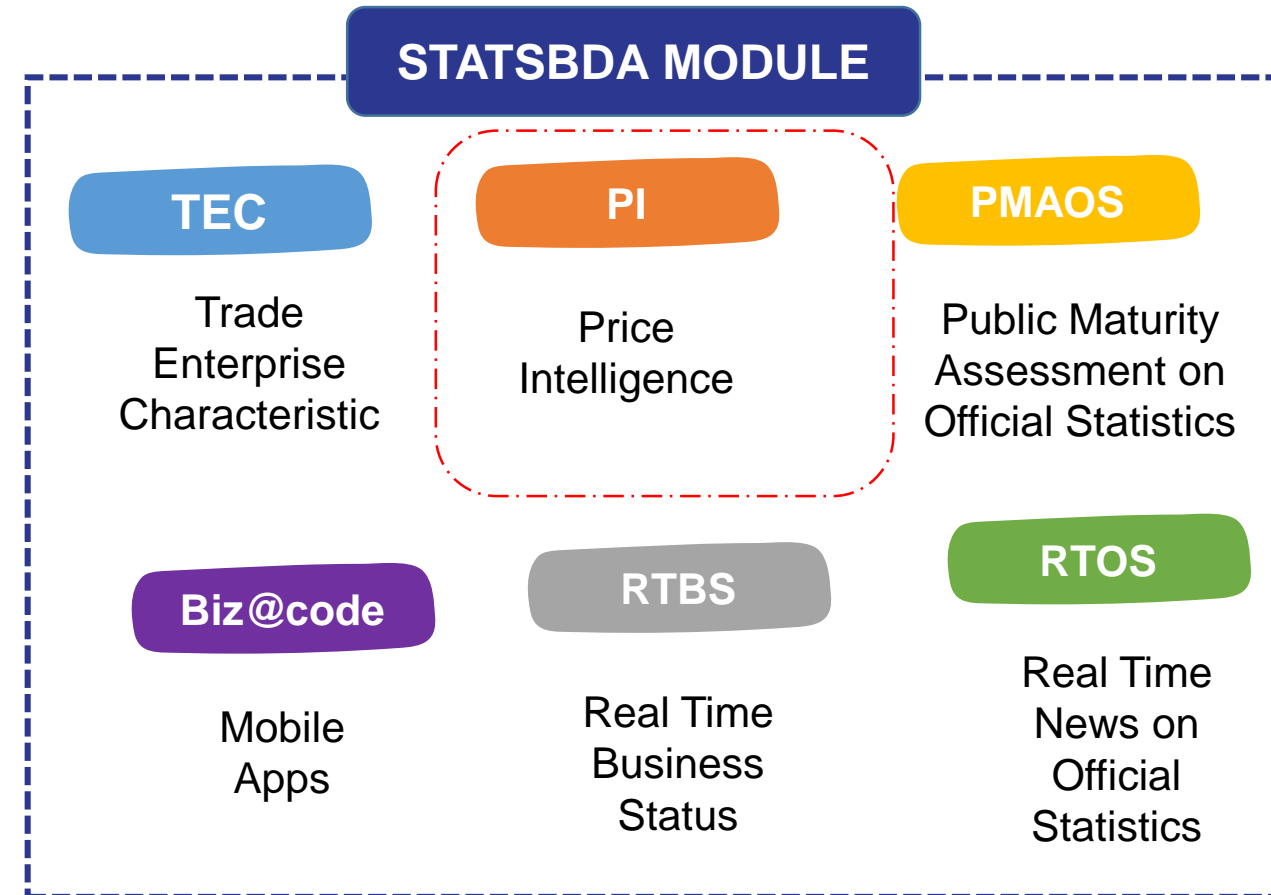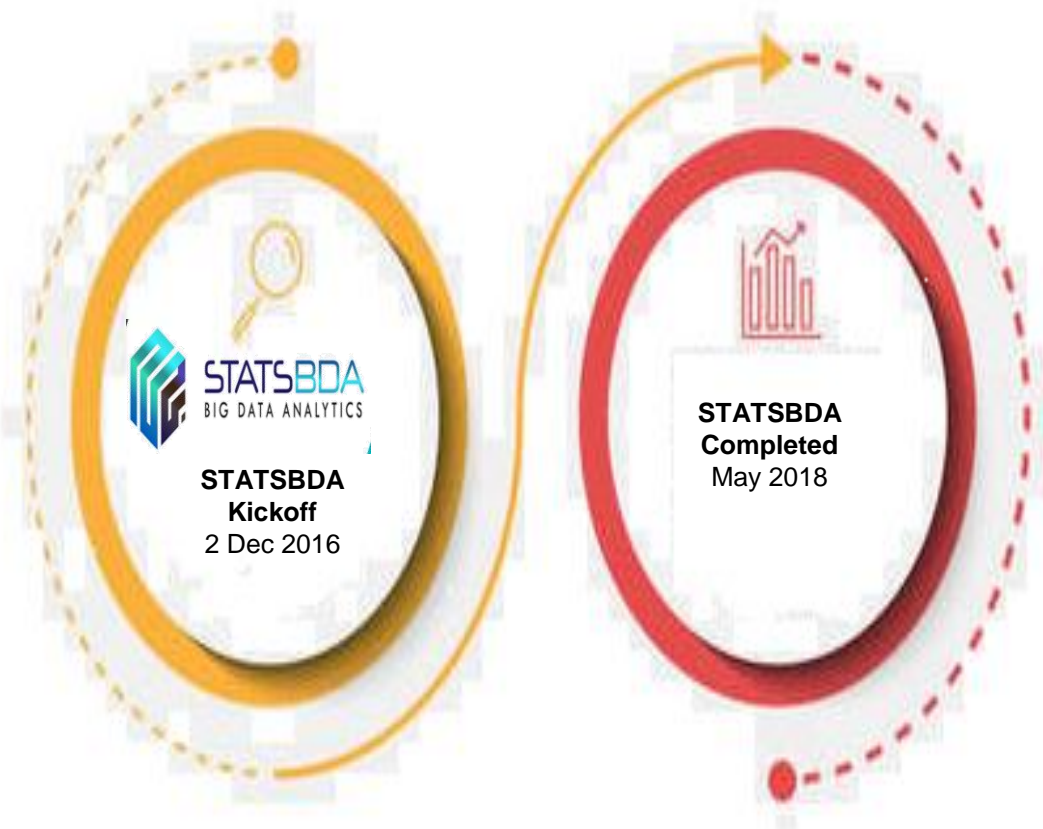
@StatsMalaysia
@MyCensus2020
www.dosm.gov.my
www.mycensus.gov.my

The Department of Statistics Malaysia (DOSM) has initiated the implementation of BDA under the project of Statistics Big Data Analytics (STATSBDA).

**STATSBDA Kickoff**
2 Dec 2016

**STATSBDA Completed**
May 2018

## STATSBDA MODULE

**TEC**
Trade Enterprise Characteristic

**PI**
Price Intelligence

**PMAOS**
Public Maturity Assessment on Official Statistics

**Biz@code**
Mobile Apps

**RTBS**
Real Time Business Status

**RTOS**
Real Time News on Official Statistics

To reduce
respondents burden

To use as
supplements for
existing data in
production of
certain statistics

To allow high
accuracy

**Objective**

To modernise data
collection

To produce new
statistical indicators

## Price Intelligence

- leveraging the capability of Big Data in collecting large data from various sources and transform them into better structure

- different prices of the same good can be obtained through various online retailer websites, providing a modernized price data collection

- Transform from unstructured data into structured data to perform analysis



web crawler: crawler → visit all links → build list → indexing → store in database

web scraping: website → scraper → xml, sql, excel data

## What is web crawler?

process of repetitively finding and fetching hyperlinks starting from a list of starting URLs.

## What is web scraping?

Web-scraping is automatically retrieving and processing information from websites

# PRICE INTELLIGENCE OBJECTIVE

## Objective

**01**

to give better insight in consumer price analysis and monitoring

**02**

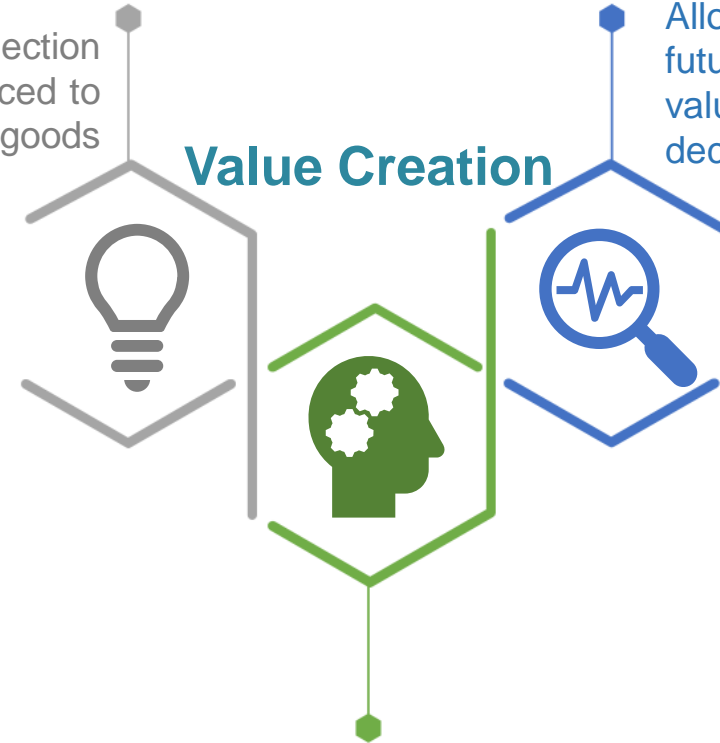to establish new price basket analysis which will be used as value added to the current Consumer Price Index

**STATSBDA**
BIG DATA ANALYTICS

PI Visualization

Data Management

**New Data Collection Methodology**
New data collection process introduced to cover price of goods selling online
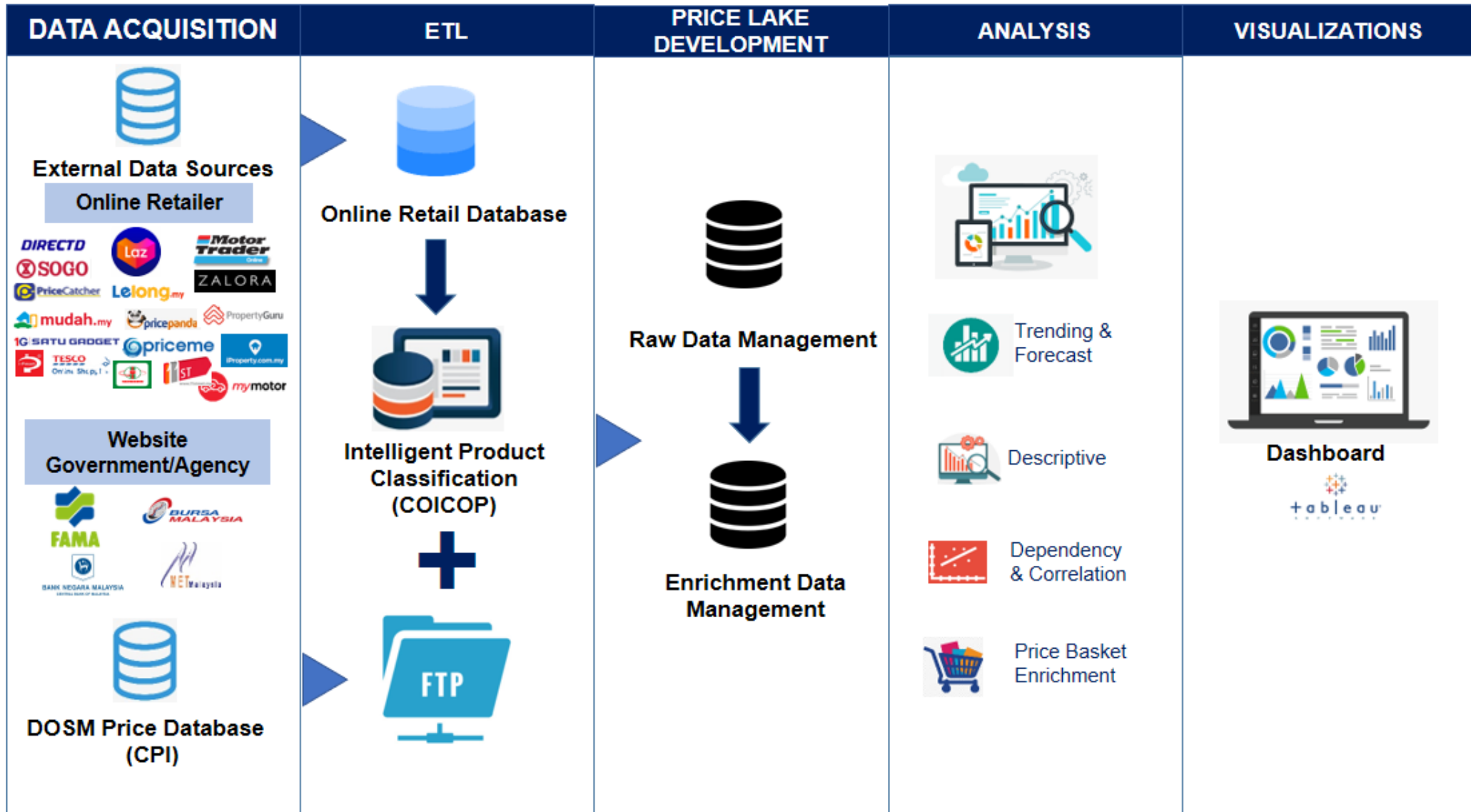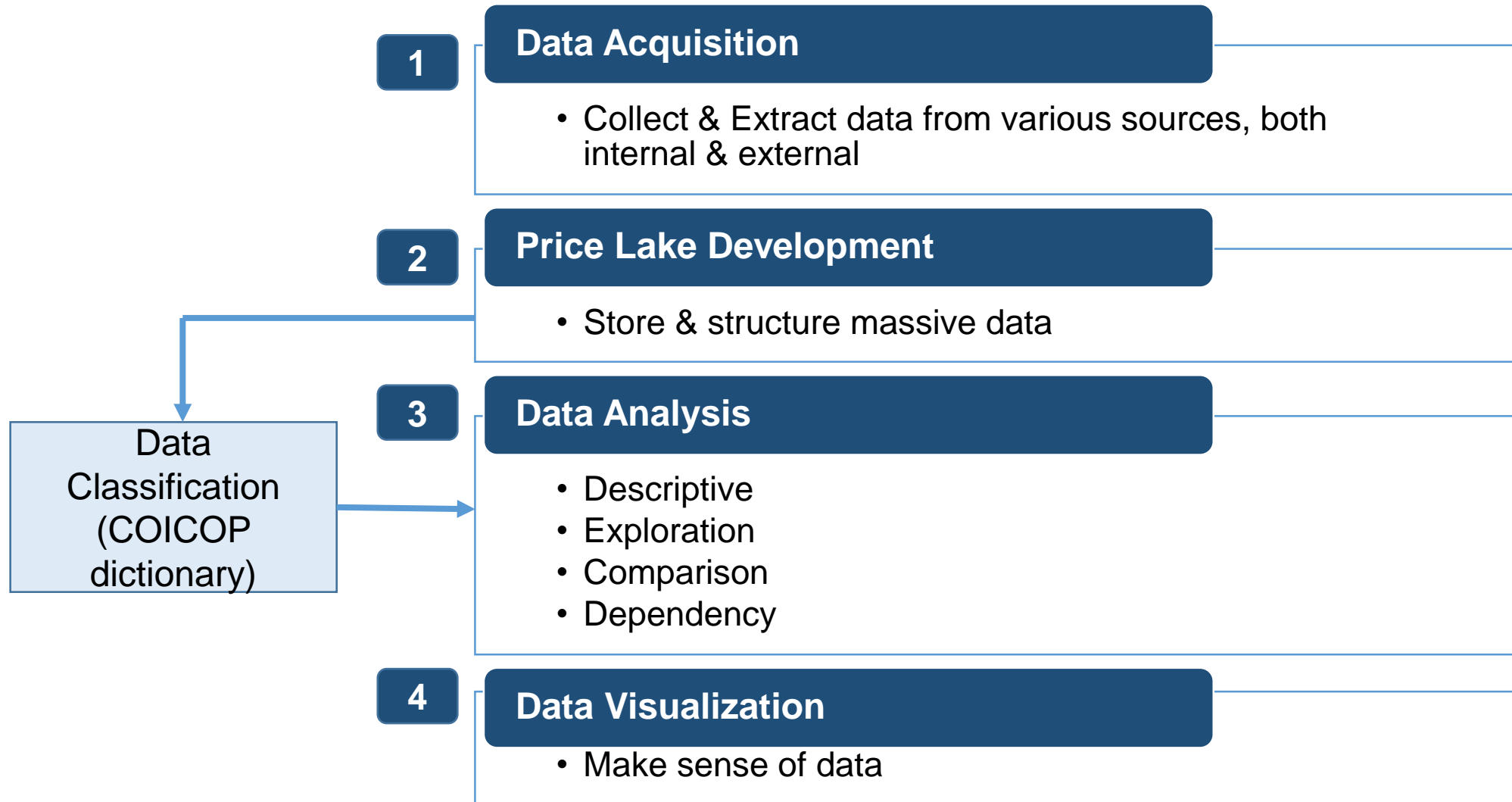
**Value Creation**

**Holistic View in Online and Offline Prices**
Allow monitoring and forecast future price trend and as valuable input for price control decisions by government

**Transform of Business Process**
New data collected enable to create holistic landscape of CPI monitoring process

# MODULE IN PRICE INTELLIGENCE

**1** **Data Acquisition**

- Collect & Extract data from various sources, both internal & external

**2** **Price Lake Development**

- Store & structure massive data

Data Classification (COICOP dictionary)

**3** **Data Analysis**

- Descriptive
- Exploration
- Comparison
- Dependency

**4** **Data Visualization**

- Make sense of data

Note : COICOP → Classification of Individual Consumption According to Purpose

## PI Dictionary



Pending/Accepted refers to the matching result at the 4 Digit COICOP level

Partially/exact refers to the matching result with MCOICOP

11.3% exact match
88.7% partially match

The word most often appear on the item description

*Information based on 231 item specifications updates in PI Dictionary

## PI Dictionary – Text matching scoring

Specification that meet requirement of items in CPI basket of goods with high scoring value will return exact match, mapped to 7 digits item specification MCOICOP code

| Item Specification | 6D Code | Item Code | Score | Match Type |
|---|---|---|---|---|
| (i) huggies dry diapers m 6-11 kg 72 pieces each | 121321 | 1213212 | 75.14 | EXACT |
| (i) drypers wee wee dry disposable diapers m 6-11kg 74pcs each | 121321 | 1213211 | 76.59 | EXACT |
| (i) huggies dry diapers m 6-11 kg 72 pieces each | 121321 | 1213212 | 75.14 | EXACT |
| (i) huggies dry diapers m 6-11 kg 72 pieces each | 121321 | 1213212 | 75.14 | EXACT |
| (i) drypers wee wee dry disposable diapers m 6-11kg 74pcs each | 121321 | 1213211 | 76.59 | EXACT |
| (i) drypers wee wee dry disposable diapers m 6-11kg 74pcs each | 121321 | 1213211 | 76.59 | EXACT |

Meanwhile, broad specification items has only 6 digits item code to be matched with online item specification. Items is considered matching as long as the item is the same regardless of brand name, units etc

| Item Specification | 6D Code | Item Code | Score | Match Type |
|---|---|---|---|---|
| (i) drypers wee wee dry disposable diapers xxl 15+kg 40pcs each | 121321 | | 67.12 | PARTIAL |
| (i) huggies dry pants diapers l 9-14kg 50pcs each | 121321 | | 64.44 | PARTIAL |
| (i) tena value diapers, medium 8 pack x 12s | 121321 | | 50.86 | PARTIAL |

## PI Dictionary



PI Dictionary has been build in order to get the best match to MCOICOP 6/7 digits based on description of the items using text matching scoring

Dictionary will be updated regularly based on basket of goods in CPI
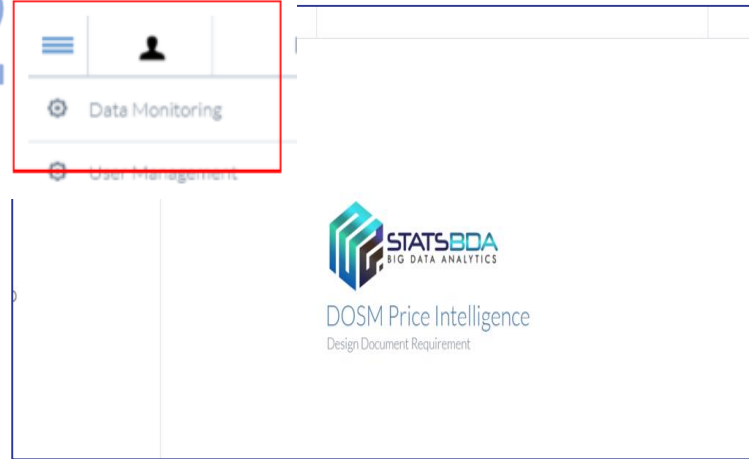
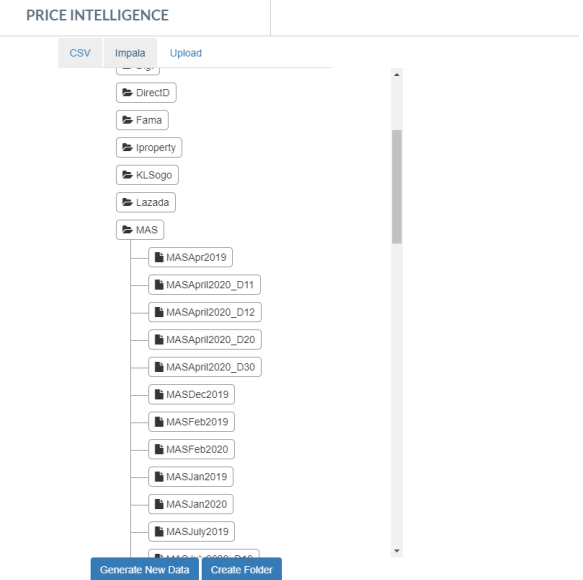COICOP : Classification of Individual Consumption According to Purpose
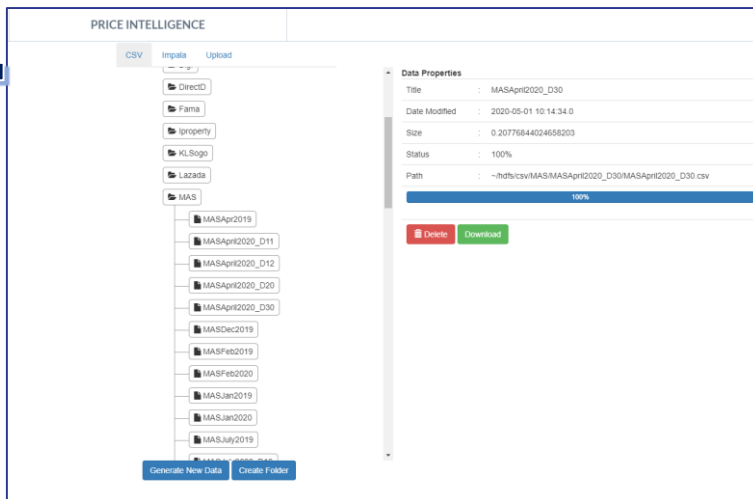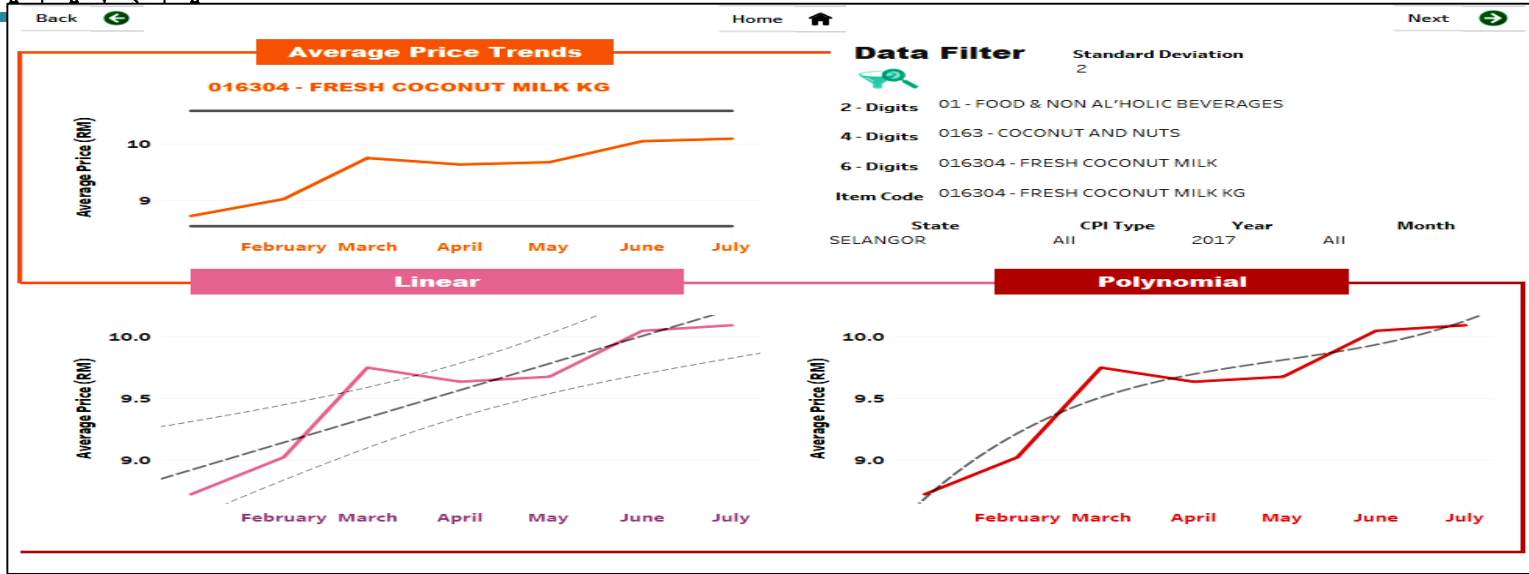
# DATA SET GENERATOR



**1**



**2**

**3**



**4**



**5**



## Download Data

1. Click Data Monitoring on the top left of the screen
2. Click Folder Name and Click file name
3. Status of the generated data will be displayed.
   Click Download to download to local PC

# PRICE INTELLIGENCE ANALYTICS



## a) Trend
Trend analysis is used to evaluate data patterns based on linear approaches. There are various time series analysis techniques that can be used, such as ARIMA, Exponential Smoothing, Holt Winter, Linear Trend, Exponential and Level Aggregation.

## b) Descriptive
Descriptive analysis, more emphasis on data exploration summaries such as: Mean, Median, Standard Deviation , Variants, Histogram / Skewness, etc.

## c) Price Basket Enrichment

Price Basket Enrichment is the process of adding data crawled from online sources into available data from the SIHP-CPI System. This process requires DOSM support because not all prices used in CPI manual calculation are available online.





## d) *Dependency*

Dependency analysis looks for more connections between data. An example is to evaluate variables by performing factor analysis or releasing irrelevant data. Cluster analysis can be carried out to assess data relationships (correlations) and other analysis can be conducted to assess such as Factor Analysis - Cause Analysis, Correlation and Cross Category

Department of Statistics
MALAYSIA

During the pandemic Covid-19, online flight ticket prices has been used in the compilation of the CPI.

## Item specification (Route)

☐ Kedah to Kuala Lumpur
☐ Johor Bahru to Kuala Lumpur
☐ Kuantan to Kuala Lumpur
☐ Penang to Kuala Lumpur
☐ Kuala Terengganu to Kuala Lumpur

Implications of the Movement Control Order (MCO), all the price data collection at the outlet has been suspended. Data crawling has been done during the MCO for 20 main product CPI as below:

**MAC 2020**

- CAP RAMBUTAN HIJAU SST RICE 5% 10KG (RICE)
- JATI BERAS SUPER SPESIAL 10KG (RICE)
- JASMINE RICE SUPER SPECIAL 5% 10KG (RICE)
- JASMINE RICE SUPER SPECIAL TEMPATAN 5KG (RICE)
- AYAM BERSIH (PELBAGAI BAHAGIAN) (CHICKEN)
- IKAN BAWAL HITAM (FISH)
- IKAN CENCARU (FISH)
- IKAN KEMBUNG (FISH)
- KUETIAU BASAH (FLAT RICE NOODLES)
- MEE KUNING BASAH (NOODLES)
- UBI KENTANG (POTATO)
- BAWANG BESAR (ONION)
- CILI KERING KERINTING (DRIED CHILLI)
- KACANG BUNCIS (FRENCH BEAN)
- KUBIS BULAT (TEMPATAN) (CABBAGE)
- LOBAK MERAH (CARROT)
- TOMATO
- SANTAN KELAPA (FRESH COCONUT MILK)
- TELUR AYAM GRED B (HEN'S EGGS GRADE B)
- *MINYAK MASAK (PELBAGAI JENAMA) (COOKING OIL)

## HOUSE RENTAL PRICES FOR INTERNATIONAL COMPARISON PROGRAM (ICP)

45.6% housing rental data for International Comparison Program 2017 (ICP2017) submission were using online price data

### RENTAL HOUSING DATA SOURCED USED FOR ICP 2017



Online Price | Rent Survey/ICP Price Collection

**19** types of housing specification with different sizes consisting of
- Single-detached house
- Attached house (row house)
- Studio apartment
- One-bedroom apartment
- Two-bedroom apartment
- Three-bedroom apartment

❏ Need to update crawler

❏ Have to build a bunch of crawlers for different sites

❏ The structure of websites change frequently

❏ Legal issues involved

❏ Storage limitation (huge amounts of data)

❏ Access and scrape data which is publicly available and avoid trying to crawl data which is private or protected by copyrights and other laws;

❏ Always check the website's robots.txt file

**Department of Statistics MALAYSIA**

TESCO

**Here to help**
Guide price
Safe online shopping
Terms & Conditions
Privacy policy

**About**
Where we deliver
Service charge
Payment options
Tesco.com.my
Clubcard

**First time shopping**
How to shop
Registration
Book a delivery
Favourites

**Contact us**
tescohelpline@tesco.com.my
1300131313
Store locator

**Intellectual Property**

The content of the Tesco Website is protected by copyright, trade marks, database and other intellectual property rights. You may retrieve and display the content of the Tesco Website on a computer screen, store such content in electronic form on disk (but not any server or other storage device connected to a network) or print one copy of such content for your own personal, non-commercial use, provided you keep intact all and any copyright and proprietary notices. You may not otherwise reproduce, modify, copy or distribute or use for commercial purposes any of the materials or content on the Tesco Website without written permission from Tesco Website.

No licence is granted to you in these Terms and Conditions to use any of our trade marks.

https://shopee.com.my/robots.txt

```
User-Agent:Googlebot
User-Agent:Bingbot
Crawl-delay:0.1
Disallow: /cart/
Disallow: /checkout/
Disallow: /buyer/
Disallow: /user/
Disallow: /me/
Disallow: /order/
Disallow: /daily_discover/
Disallow: /mall/just-for-you/
Disallow: /mall/*-cat.
Disallow: /from_same_shop/
Disallow: /you_may_also_like/
Disallow: *-i.*/similar?from=flash_sale
Disallow: /find_similar_products/
Disallow: /top_products
Disallow: /search*searchPrefill
Disallow: /index.html         .

User-Agent:*
Crawl-delay:1
Disallow: /cart/
Disallow: /checkout/
Disallow: /buyer/
Disallow: /user/
Disallow: /me/
Disallow: /order/
Disallow: /daily_discover/
Disallow: /mall/just-for-you/
Disallow: /mall/*-cat.
Disallow: /from_same_shop/
Disallow: /you_may_also_like/
Disallow: *-i.*/similar
Disallow: /find_similar_products/
Disallow: /top_products
Disallow: /search*searchPrefill
Disallow: /index.html
```

A robots.txt file tells search engine crawlers which pages or files the crawler can or can't request from your site. This is used mainly to avoid overloading your site with requests; **it is not a mechanism for keeping a web page out of Google.** To keep a web page out of Google, you should use noindex directives, or password-protect your page.

A robots.txt file is used primarily to manage crawler traffic to your site, and *usually* to keep a file off Google, depending on the file type.

Source:
*https://developers.google.com/search/docs/advanced/robots/intro*

**1** **Forecasting prices of fish and vegetable using web scraped price micro data**

- Mazliana Mustapa, Raja Rajeswari Ponnusamy, Ho Ming Kang

**2** **Online and Offline Prices: Measuring selected home appliance's products**
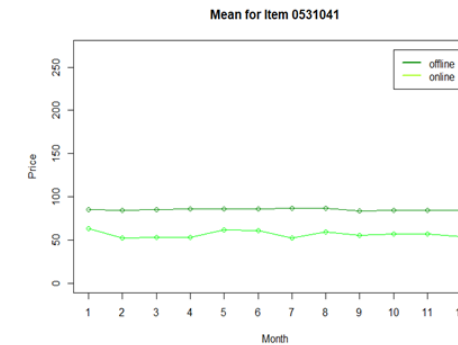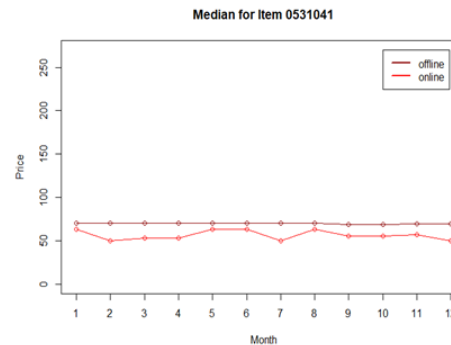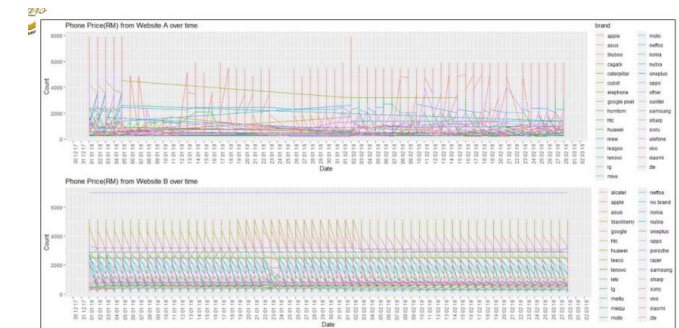
- Mohd Saiful Husain and Norsyela Muhammad Noor Mathivanan

**3** **Analysis of the Mobile Phones Prices Malaysia using Web Scraped Data**

- Nur Hurriyatul Huda Abdullah Sani



**Table 3: ARIMA best model**

| Item | Model |
|------|-------|
| Red Bream | $Y_t = 2.4205 - 0.0333Y_{t-1} + 0.6382Y_{t-2} + 0.0594 + 0.9688e_{t-1}$ |
| Selar Kuning | $Y_t = 1.4794 - 1.6648Y_{t-1} - 0.7431Y_{t-2} + e_t - 0.7924e_{t-1}$ |
| Green Spinach | $Y_t = 1.1813 + 1.7246Y_{t-1} - 0.7774Y_{t-2} + 0.1008 - 0.8126e_{t-1}$ |
| Kangkung | $Y_t = 1.0995 + 1.7293Y_{t-1} - 0.7776Y_{t-2} + 0.1023 - 0.8315e_{t-1}$ |
| Long Beans | $Y_t = 0.5460 + 0.3667Y_{t-1} + 0.0780 + 0.4322e_{t-1} + 0.2638e_{t-2}$ |
| Bawal, Cencaru, Kembong, Round Cabbage and Sawi Jepun | $Y_t = \mu + Y_{t-1}$ where $\mu$ : mean of the changes of period to period |

https://statsbda.dosm.gov.my/

**Mazliana Mustapa**
Core Team Big Data Analytics
Department of Statistics Malaysia

**Team Price Intelligence**
mazliana@dosm.gov.my
ridhuan@dosm.gov.my
noradilah.adnan@dosm.gov.my

# THANK YOU

# BANCI MALAYSIA